

The embodiments of the invention in which an exclusive property or privilege is claimed are defined as follows:

1. In a computer system comprising a cluster of node boards, each node board having at least one central processor unit (CPU) and shared memory, said node boards being interconnected into groups of node boards providing access between the central processing units (CPUs) and shared memory on different node boards, a scheduling system to schedule a job to said node boards which have resources to execute the jobs, said batch scheduling system comprising:

a topology monitoring unit for monitoring a status of the CPUs and generating status information signals indicative of the status of each group of node boards;

a job scheduling unit for receiving said status information signals and said jobs, and, scheduling the job to one group of node boards on the basis of which group of node boards have the resources required to execute the job as indicated by the status information signals.

2. The scheduling system as defined in claim 1 wherein the status information signals indicate which CPUs in each group of node boards have available resources, and, the job scheduling unit schedules jobs to groups of node boards which have resources required to execute the job.

3. The scheduling system as defined in claim 1 wherein the status information signals for each group of node boards indicate a number of CPUs available to execute jobs for each radius; and

wherein the job scheduling unit allocates the jobs to the one group of node boards on the basis of which group of node

boards have CPUs available to execute jobs of a radius required to execute the job.

4. The batch scheduling system as defined in claim 3 wherein said cluster of node boards are located on separate hosts; and wherein the topology monitoring unit monitors the status of the CPUs in each host and generates status information signals regarding groups of node boards in each host.

5. The batch scheduling system as defined in claim 4 wherein the status information signals include, for each host, a number of CPUs which are available for each radius; and wherein the scheduling unit maps the job to a selected host having a maximum number of CPUs available at a radius corresponding to the required radius for the job.

6. The batch scheduling system as defined in claim 5 further comprising, for each host, a job execution unit for receiving jobs which have been scheduled to the selected host by the job scheduling unit, and, allocating the jobs to the selected group of node boards; and

wherein the job execution unit communicates with the topology monitoring unit to allocate the jobs to the group of node boards which the topology monitoring unit has determined have the resources required to execute the job.

7. The batch scheduling system as defined in claim 1 wherein the scheduler comprises a standard scheduler for allocating jobs to the selected group of node boards and an external scheduler for receiving the status information signals from the topology monitoring unit and selecting the selected group of node boards based on the status of the information signals.

8. The batch scheduling system as defined in claim 3 wherein if the job scheduling unit cannot locate a group of node boards which have the resources required to execute the job, the job scheduling unit delays allocation of the job until the status information signals indicate the resources required to execute the job are available.

9. The batch scheduling system as defined in claim 3 wherein the access between the central processing units (CPUs) and shared memory on different node boards is non-uniform.

10. In a computer system comprising resources physically located in more than one module, said resources including a plurality of processors being interconnected by a number of interconnections in a physical topology providing non-uniform access to other resources of said computer system, a method of scheduling a job to said resources, said method comprising the steps of:

- (a) periodically assessing a status of the resources and sending status information signals indicative of the status of the resources to a job scheduling unit;
- (b) assessing, at the job scheduling unit, the resources required to execute a job;
- (c) comparing, at the job scheduling unit, the resources required to execute the job and resources available based on the status information signals; and
- (d) scheduling the job to the resources which are available to execute the job as based on the status information signals and the physical topology, and the resources required to execute the job.

11. The method as defined in claim 10 further comprising the sub-steps of :

(a)(i) periodically assessing the status of resources in each module and sending status information signals indicative of the status of the resources in each module to the job scheduling unit;

(c)(i) comparing the available resources in each module to the resources required to execute the job; and

(d)(i) scheduling the job to the module having the most resources available to execute the job.

12. The method as defined in claim 10 further comprising the sub-steps of:

(a)(i) for each module, periodically assessing the status of the resources by assessing the status of each processor in each module and sending to the job scheduling unit module status information for each module indicative of a number of available processors at each radius in the module;

(b)(i) assessing, at the job scheduling unit, the requirements necessary to execute the job by determining the number of processors of a required radius required to execute the job;

(c)(i) comparing the resources required to execute the job and the resources available by comparing the number of processors of the required radius to execute the job and the number of available processors of the required radius at each module based on the module information status signals; and

(d)(i) scheduling the job to the module which has a largest number of available processors at the required radius based on the module status information signals and the physical topology.

13. In a computer system comprising resources including a plurality of processors, said processors being interconnected by a number of interconnections in a physical topology providing non-uniform access to other resources of said computer system, a scheduling system to schedule jobs to said resources, said scheduling system comprising:

a topology monitoring unit for monitoring a status of the processors and generating status information signals indicative of the status of said processors;

a job scheduling unit for receiving said status information signals and said jobs, and, scheduling the jobs to groups of processors on the basis of the physical topology and the status information signals.

14. The scheduling system as defined in claim 13 wherein the job scheduling unit schedules the jobs based on predetermined criteria, said predetermined criteria including the expected delay to transfer information amongst the group of processors based on the physical topology and the status information signals.

15. The scheduling system as defined in claim 14 wherein the predetermined criteria include a radius of the group of processors to execute the job.

16. The scheduling system as defined in claim 15 wherein the predetermined criteria further include the number of connections in the physical topology within the group of processors, availability of memory associated with the group of processors and availability of other processors connected to the group of processors.

17. The scheduling system as defined in claim 13 wherein the plurality of processors are physically located in separate modules. Wherein the topology monitoring unit comprises topology daemons associated with each module for monitoring a status of the processors physically located in the associated module and generating module status information signals indicative of the status of the processors in the associated module, wherein the job scheduling unit receives the module status information signals from all of the topology daemons and allocates the jobs to a group of processors in one of the modules on the basis of the physical topology of the processors in the modules and the module status information signals from all of the modules.

18. The scheduling system as defined in claim 17 wherein the modules are interconnected by a META router operating on a network;

wherein the jobs and the module status information signals are communicated through the META router and network.

19. The scheduling system as defined in claim 18 wherein the network comprises an Internet.

20. The scheduling system is defined in claim 17 wherein the module status information signals indicate a number of available processors for each radius; and wherein the job scheduling unit schedules the job to a module having available processor of a radius required to execute the job.

21. The scheduling system as defined in claim 20 wherein the

scheduling unit schedules jobs to the module having a greatest number of available processors of a radius required to execute the job.